



041 radiology report generation (RRG) has received much at-  
042 tention [36, 43, 44, 51].

043 Traditional class imbalance is a prevalent issue in vari-  
044 ous medical datasets. Due to factors such as disease inci-  
045 dence rates and diagnosis difficulty, the data scale of dif-  
046 ferent classes in medical datasets often vary significantly,  
047 making it challenging for the model to learn the features  
048 of minority classes. However, it is noteworthy that unlike  
049 general classification tasks, RRG is more complex because  
050 its labels are long texts containing extensive information,  
051 and although the classes are fixed, the content of the reports  
052 varies. Consequently, we are prompted to raise two pre-  
053 viously overlooked questions: **Q1: What are the impacts of**  
054 **traditional class imbalance on RRG models? Q2: Is there a**  
055 **better metric to more effectively measure the imbalance in**  
056 **RRG data?**

057 To address **Q1**, we selected four classical and state-of-  
058 the-art (SOTA) RRG methods [7, 12, 16, 47] and conducted  
059 experiments on the currently largest publicly available RRG  
060 dataset, MIMIC-CXR [21, 22]. We divided the test set into  
061 13 classes based on the official tags and evaluated the qual-  
062 ity of the generated reports (BLEU-4 score) for each class.  
063 We arranged the classes in descending order of their sam-  
064 ple scale (gray bars) and observed the relationship with the  
065 quality of the generated reports, as shown in Figure 1(a).  
066 Intuitively, according to the traditional class imbalance as-  
067 sumption, the fewer training samples a class has, the worse  
068 its corresponding generated reports will be. In other words,  
069 the two should be almost proportional. However, as shown  
070 in the figure, this assumption does not hold significantly for  
071 the RRG task.

072 Based on the above finding, we propose a method  
073 of **Radiology Report Generation by Curriculum Learning**  
074 **(RRGCL)** to address **Q2**. Inspired by curriculum learn-  
075 ing (CL) [3], we conceive a novel learning difficulty met-  
076 ric to reevaluate imbalance in RRG data. This metric com-  
077 prises two components: (1) visual neighborhood seman-  
078 tic gap (VNS-Gap); (2) image-text cross-modal alignment  
079 score (ITAS). We rearrange the classes in descending order  
080 according to the proposed metric and again observe their  
081 relationship with the quality of the generated reports, with  
082 the results shown in Figure 1(b). As evidenced in the fig-  
083 ure, the proposed metric demonstrates a stronger positive  
084 correlation with report generation quality, indicating its su-  
085 perior efficacy in measuring imbalance within RRG data.  
086 Furthermore, we design a novel sample scheduling function  
087 based on the proposed imbalance metric. It dynamically  
088 adjusts the pace of introducing data of different difficulties  
089 into training subset, thereby helping the model gradually  
090 adapt to the inherent imbalance in RRG data. Extensive  
091 experiments demonstrate that RRGCL can effectively help  
092 existing RRG models further improve the quality of gener-  
093 ated reports. Overall, the main contributions of this paper

are summarized as follows:

- We discover that the traditional class imbalance based on  
sample quantity has limitations in RRG data, and propose  
a novel learning difficulty metric that can more effectively  
measure imbalance in RRG data.
- To address the redefined imbalance issue, we propose a  
simple yet efficient sample scheduling strategy for RRG  
models.
- Our proposed RRGCL method is a model-agnostic and  
plug-and-play training strategy. Extensive experiments  
demonstrate that existing RRG methods achieve perfor-  
mance improvements after applying RRGCL.

## 2. Related Work

### 2.1. Radiology Report Generation

Automatic radiology report generation (RRG) originates  
from image captioning task, but it is more complex and  
challenging. Due to pressing clinical needs, research in  
RRG has been steadily increasing in recent years [44]. Cur-  
rent RRG studies can be broadly categorized into three  
main classes: CNN-RNN-based method, Transformer-  
based method and retrieval-based method.

CNN-RNN-based models [13, 35, 38, 47, 52, 56, 61]  
represent the most classic encoder-decoder frameworks for  
RRG. In this paradigm, CNN network or its variants (e.g.,  
ResNet, DenseNet) are typically employed to extract vi-  
sual features, which are then decoded by RNN network  
or its variants (e.g., LSTM, GRU) to generate the re-  
port word by word. Such RRG methods are adapted from  
early image captioning tasks, laid a crucial foundation  
for subsequent research, and continue to be influential to-  
day. However, they often exhibit limitations in modeling  
long-range dependencies and tend to overlook subtle le-  
sions in images. The introduction of Transformer [48] ef-  
fectively mitigates this issue. Transformer-based methods  
[5, 7, 8, 12, 16, 24, 27, 31, 34, 37, 49, 53] represent the do-  
minant paradigm in current RRG research. These approaches  
typically employ a Transformer as the decoder to generate  
long-text reports, while the image encoder is usually a vi-  
sual Transformer (ViT) [9] or a hybrid CNN-Transformer  
architecture for visual feature extraction. Thanks to the at-  
tention mechanism, such methods are often better at mod-  
eling long-range dependencies, thereby generating higher-  
quality and more human-like reports. Due to similar pat-  
terns among radiology reports, some retrieval-based meth-  
ods [10, 19, 46, 58, 62] have been introduced. These models  
first retrieve similar reports or sentences from a template or  
report corpus, and then use this as reference knowledge to  
generate the final report. These methods are heavily depen-  
dent on large-scale, specialized databases and exhibit lim-  
ited generalization capability.

Our proposed RRGCL method falls outside the afore-

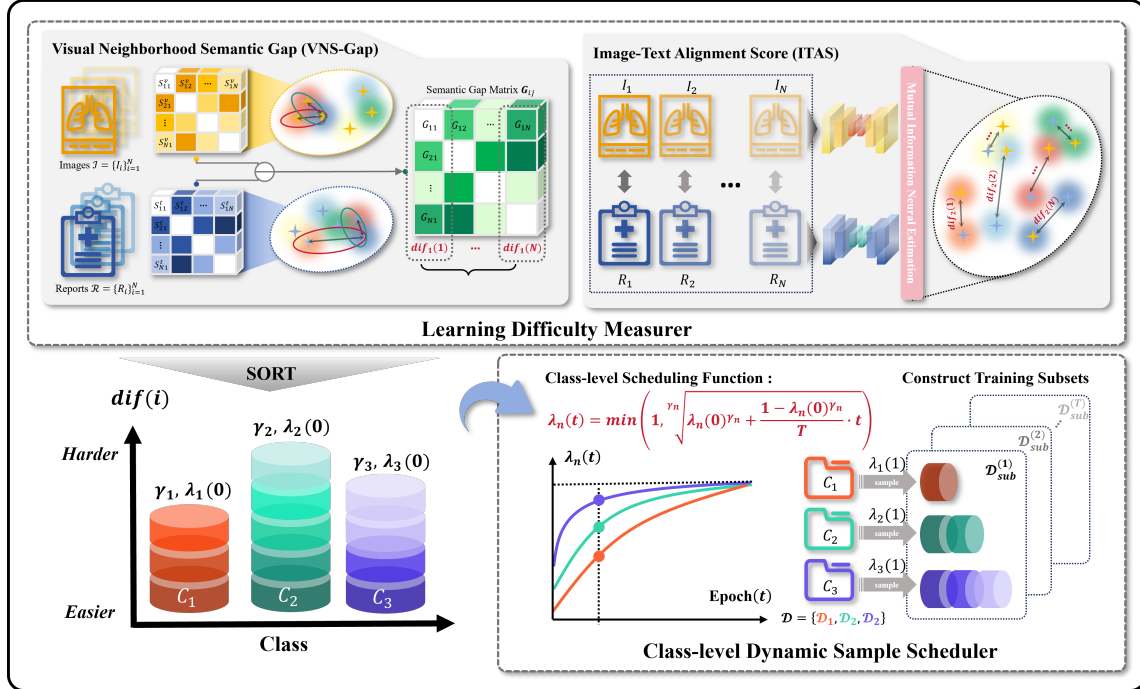


Figure 2. The overall framework of the proposed RRGCL primarily consists of two components: a Learning Difficulty Measurer and a Class-level Dynamic Sample Scheduler. The measurer evaluates the learning difficulty  $dif(i)$  of each sample by calculating VNS-Gap  $dif_1(i)$  and ITAS  $dif_2(i)$ . Based on  $dif(i)$ , all training samples are sorted from easy to hard within each class. The scheduler dynamically samples from each class in sequence according to our designed class-level scheduling function  $\lambda_n(t)$ . Simpler classes are incorporated into the training subset at a faster rate during the early training stages, while more difficult classes are introduced at a more gradual pace. Class difficulty scores reflects the underlying imbalance in RRG data from a novel perspective, and the proposed scheduling strategy effectively addresses this imbalance.

145 mentioned classes, because it cannot generate reports by it-  
 146 self. Instead, it is a model-agnostic training strategy that  
 147 can be applied to various existing RRG models to further  
 148 enhance their performance.

## 149 2.2. Curriculum Learning

150 The core idea of curriculum learning (CL) [3] is to mimic  
 151 the human cognitive process of progressing from easy to  
 152 difficult. In the early stages of training, the model primar-  
 153 ily encounters the simplest, cleanest, and most representa-  
 154 tive samples in the dataset, and then gradually incorporates  
 155 more complex data. This strategy provides the model with  
 156 a better initial optimization path, leading to more stable and  
 157 efficient training [45, 50].

158 CL has been proven effective in various downstream  
 159 tasks. For example, in the field of computer vision (CV),  
 160 it has been applied successfully to medical image classifica-  
 161 tion [20, 26], face recognition [15], object detection [59],  
 162 etc. In the field of natural language processing (NLP), it has  
 163 shown effectiveness in machine translation [11], text emo-  
 164 tion recognition [25, 57], text generation [6], etc. Addi-  
 165 tionally, it has also been applied to some cross-modal tasks  
 166 like image caption [55] and video caption [28]. To the best

of our knowledge, CMCL [32] is currently the only CL-  
 based method for RRG. CMCL defines two training diffi-  
 culty evaluation metrics for images and text respectively.  
 During training, it sorts all training data into four batches  
 based on these metrics and dynamically selects the most  
 appropriate batch for training according to the model’s cur-  
 rent learning competence. Our proposed method does not  
 evaluate images and text in isolation. Instead, it calculates  
 the learning difficulty of each sample by leveraging inter-  
 sample divergence and image-text relationship.

## 3. Method

RRGCL mainly consists of two parts: a learning difficulty  
 measurer and a dynamic sample scheduler. The learning  
 difficulty measurer quantitatively evaluates how challeng-  
 ing it is to learn from a class, thereby revealing the un-  
 derlying imbalance in the RRG data. The dynamic sam-  
 ple scheduler determines how to progressively increase the  
 learning difficulty during training, which we utilize to miti-  
 gate the imbalance issue. The overall framework of RRGCL  
 is shown in Figure 2.

### 3.1. Definition

Let  $\mathcal{D} = \{I_i, R_i, y_i\}_{i=1}^N$  denote the training dataset including  $N$  samples, where  $I_n \in \mathcal{I}$  represents the  $i$ -th radiology image and  $R_i \in \mathcal{R}$  represents the corresponding  $i$ -th report, and  $y_i$  is the label including  $c$  classes. The objective of RRG is to learn a mapping function  $\mathcal{G}_\theta : \mathcal{I} \rightarrow \mathcal{R}$  that generates reports from input images.

## 3.2. Learning Difficulty Measurer

### 3.2.1. What Does Difficult RRG Data Look Like?

Before elaborating on our method, it is necessary to consider this question: *What kind of RRG data might be difficult for the model to learn?* We answer this question from two perspectives. (1) Visual differences between medical images are often subtle, yet the descriptions of these minor variations in reports can be substantially distinct. For example, a slight increase in the density of a region might be described as new-onset inflammation in the report. Therefore, if an RRG sample exhibits visual similarity to other cases but is associated with a significantly divergent report, it is likely difficult to learn, as such samples create confusion for the model. (2) RRG is fundamentally a cross-modal transfer task from visual to textual data, where the degree of inter-modal alignment directly influences performance. Thus, samples with stronger image-text alignment are generally easier to learn. Based on these considerations, we design two specialized metrics to quantify the learning difficulty of RRG data: Visual Neighborhood Semantic Gap (VNS-Gap) and Image-Text Alignment Score (ITAS).

### 3.2.2. Visual Neighborhood Semantic Gap

The visual neighborhood semantic gap (VNS-Gap) quantifies the learning difficulty by measuring the discrepancy between visual similarity and semantic similarity across samples. Let  $\phi : \mathcal{I} \rightarrow \mathbb{R}^d$  be a pre-trained visual encoder, and the visual similarity between two images  $I_i$  and  $I_j$  is computed as follows:

$$\mathbf{S}_{ij}^v = \frac{\phi(I_i)^\top \phi(I_j)}{\|\phi(I_i)\|_2 \cdot \|\phi(I_j)\|_2} \quad (1)$$

where  $\mathbf{S}_{ij}^v \in \mathbb{R}^{N \times N}$  represents the visual similarity matrix. Similarly, let  $\psi : \mathcal{R} \rightarrow \mathbb{R}^d$  be a pre-trained textual encoder, and the textual similarity between two reports  $R_i$  and  $R_j$  is computed as follows:

$$\mathbf{S}_{ij}^t = \frac{\psi(R_i)^\top \psi(R_j)}{\|\psi(R_i)\|_2 \cdot \|\psi(R_j)\|_2} \quad (2)$$

where  $\mathbf{S}_{ij}^t \in \mathbb{R}^{N \times N}$  represents the textual similarity matrix. To capture the absolute discrepancy between visual and textual similarities, a semantic gap matrix  $\mathbf{G}_{ij} \in \mathbb{R}^{N \times N}$  is calculated as follows:

$$\mathbf{G}_{ij} = |\mathbf{S}_{ij}^v - \mathbf{S}_{ij}^t| \quad (3)$$

Finally, the VNS-Gap of the  $i$ -th sample can be easily measured as follows:

$$dif_1(i) = \sum_{j=1}^N \mathbf{G}_{ij} \quad (4)$$

### 3.2.3. Image-Text Alignment Score

Mutual Information (MI) is an efficient statistical measure for calculating the image-text alignment score (ITAS) between radiology images and corresponding reports. The MI between image  $I$  and report  $R$  is defined as follows:

$$M(I; R) = D_{KL}(\mathbb{P}_{IR} \| \mathbb{P}_I \otimes \mathbb{P}_R) \quad (5)$$

where  $\mathbb{P}_{IR}$  denotes the joint distribution of image-text pairs,  $\mathbb{P}_I \otimes \mathbb{P}_R$  represents the product of marginal distributions, and  $D_{KL}$  is the Kullback-Leibler (KL) divergence. However, estimating MI between high-dimensional continuous variables with limited data is quite challenging [29], so we employ Mutual Information Neural Estimation (MINE) [18] to circumvent the issue of the dimensionality disaster. MINE leverages the Donsker-Varadhan (DV) dual representation of the KL-divergence, which provides a lower bound for MI as follows:

$$M(I; R) \geq \sup_{T \in \mathcal{F}} [\mathbb{E}_{\mathbb{P}_{IR}} [T(I, R)] - \log \left( \mathbb{E}_{\mathbb{P}_I \otimes \mathbb{P}_R} \left[ e^{T(I, R)} \right] \right)] \quad (6)$$

where  $T : \mathcal{I} \times \mathcal{R} \rightarrow \mathbb{R}$  is a function in the family  $\mathcal{F}$  of bounded measurable functions. We parameterize the function  $T$  using a neural network  $T_\theta$  with parameters  $\theta$ , transforming the estimation into an optimization problem as follows:

$$\hat{M}_\theta(I; R) = \mathbb{E}_{\mathbb{P}_{IR}} [T_\theta(I, R)] - \log \left( \mathbb{E}_{\mathbb{P}_I \otimes \mathbb{P}_R} \left[ e^{T_\theta(I, R)} \right] \right) \quad (7)$$

The optimal parameters are obtained by solving:

$$\theta^* = \arg \max_{\theta} \hat{M}_\theta(I; R) \quad (8)$$

After training, we obtain an image-text alignment scorer  $T_{\theta^*}$ , which evaluates  $i$ -th sample  $(I_i, R_i)$  as follows:

$$dif_2(i) = T_{\theta^*}(I_i, R_i) \quad (9)$$

### 3.2.4. Overall Learning Difficulty

The overall learning difficulty score for  $i$ -th sample is obtained by combining both VNS-Gap and ITAS after normalization:

$$dif(i) = \alpha \cdot \frac{dif_1(i) - \mu_1}{\sigma_1} + \beta \cdot \frac{dif_2(i) - \mu_2}{\sigma_2} \quad (10)$$

where  $(\mu_1, \sigma_1)$  and  $(\mu_2, \sigma_2)$  are the mean and standard deviation of  $dif_1(i)_{i=1}^N$  and  $dif_2(i)_{i=1}^N$ , respectively.  $\alpha, \beta \in \mathbb{R}^+$  are weighting coefficients, which will be discussed in Section 4.6.

### 273 3.3. Class-level Dynamic Sample Scheduler

274 Most sample scheduler is to sort all training data from easy  
275 to hard based on learning difficulty and then schedule the  
276 data sequentially according to a certain rule. This paradigm  
277 may result in some classes with fewer samples rarely or  
278 never appearing in the newly added training subsets during  
279 the later training stages, which is likely to exacerbate data  
280 imbalance. To end this, we propose a class-level dynamic  
281 sample scheduling strategy, as shown in Algorithm 1.

---

#### Algorithm 1 Sample Scheduling Strategy of RRGCL

---

**Input:** The training set  $\mathcal{D}$ , the number of training epochs  $T$ , the training difficulty  $diff(i)$  for each sample.

**Output:** A RRG model  $\mathcal{G}_\theta$ .

- 1: Divide  $\mathcal{D}$  into  $c$ -class subsets, resulting in  $\mathcal{D} = \{\mathcal{D}_n\}_{n=1}^c$ ;
  - 2: Sort samples in  $\mathcal{D}_n$  based on the sample difficulty and calculate class difficulty  $diff_C^{(n)}$  based on  $diff(i)$ ;
  - 3: Calculate “scheduling exponential”  $\gamma_n$  and initial scheduling rate  $\lambda_n(0)$ ;
  - 4: Calculate scheduling function  $\lambda_n(t)$  by Eq.11;
  - 5:  $\mathcal{D}_{sub}^{(t)} = \emptyset$ ;
  - 6: **for**  $t = 0$  to  $T$  **do**
  - 7:   **for**  $n = 1$  to  $c$  **do**
  - 8:     Extract the top  $\lambda_n(t)$  proportion of training samples from each class;
  - 9:      $\mathcal{D}_{sub}^{(t)} \leftarrow \mathcal{D}_{sub}^{(t+1)}$ ;
  - 10:   **end for**
  - 11:   Train  $\mathcal{G}_\theta$  on  $\mathcal{D}_{sub}^{(t)}$ ;
  - 12: **end for**
- 

#### 282 3.3.1. Designing Class-level Scheduling Function

283 We divide the training set  $\mathcal{D}$  into  $c$  subsets according to the  
284 class label  $y_i$ , and then arrange the samples within each sub-  
285 set in ascending order of learning difficulty  $diff(i)$ , resulting  
286 in  $\mathcal{D} = \{\mathcal{D}_n\}_{n=1}^c$ . A class-level learning difficulty metric  
287  $diff_C^{(n)}$  is introduced to measure the difficulty of each class  
288 by averaging the difficulties of all samples within that class.

289 The root function [42] has been proven to be an effective  
290 function for sample scheduling. Our proposed scheduling  
291 strategy assigns each class a root function  $\lambda_n(t)$  based on its  
292 class-level difficulty, thereby controlling the sampling rate  
293 from different classes. The specific is as follows:

$$294 \lambda_n(t) = \min \left( 1, \sqrt[\gamma_n]{\lambda_n(0)^{\gamma_n} + \frac{1 - \lambda_n(0)^{\gamma_n}}{T} \cdot t} \right) \quad (11)$$

295 where  $T$  represents the total number of training epochs, the  
296 variable  $t \in [1, T]$  represents the current training epoch,  
297 and the value range of  $\lambda_n(t)$  is  $(0, 1]$ .  $\gamma_n \in [1, c + 1]$  is a  
298 “scheduling exponential” we defined, which is a function of

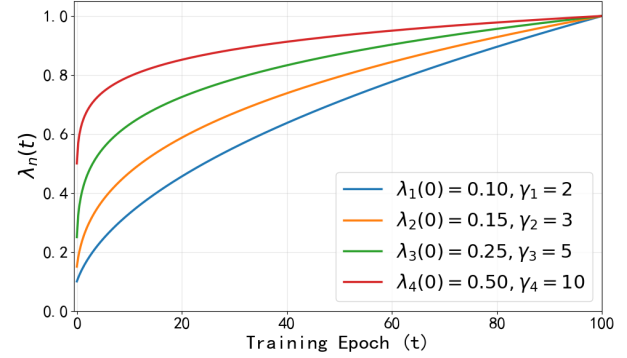


Figure 3. An example of scheduling function  $\lambda_n(t)$  over training epochs for different  $\gamma_n$  and  $\lambda_n(0)$ , where  $T = 100$ .

the class-level difficulty  $diff_C^{(n)}$ , as shown below: 299

$$\gamma_n \left( diff_C^{(n)} \right) = 1 + (1 - diff_C^{(n)}) \cdot c \quad (12) \quad 300$$

This is a simple yet efficient method that ensures simpler 301  
classes are assigned larger “scheduling exponential”.  $\lambda_n(0)$  302  
is the initial scheduling rate for the  $n$ -th class can be set 303  
according to  $\gamma_n$  as follows: 304

$$\lambda_n(0) = \frac{\gamma_n}{\sum_{i=1}^c \gamma_i} \quad (13) \quad 305$$

#### 3.3.2. Constructing Training Subsets 306

307 For each training epoch  $t$ , the scheduler constructs a training  
308 subset  $\mathcal{D}_{sub}^{(t)}$  by sequentially selecting the top  $\lambda_n(t)$  propor-  
309 tion of samples from each class. To intuitively illustrate the  
310 scheduling strategy we designed, a straightforward example  
311 with manually configured parameters is provided, as shown  
312 in Figure 3. It can be observed that simpler classes ( $\gamma_n$  is  
313 bigger) are scheduled in a more “aggressive” manner, with  
314 a faster scheduling rate in the early training stages, while  
315 difficult classes ( $\gamma_n$  is smaller) are scheduled more gradu-  
316 ally, progressing at a steady rate throughout the training  
317 process. At the  $T$ -th epoch, all samples join the training.  
318 By constructing training subsets in this manner, the model  
319 is exposed to simpler samples in the early training stages.  
320 In the later stages, the newly added samples still encompass  
321 data from all classes. Furthermore, due to the characteristics  
322 of the root function, the increase in difficulty throughout the  
323 training process is more gradual, ensuring training stability.

## 4. Experiment 324

### 4.1. Datasets 325

326 All experiments are conducted on the MIMIC-CXR dataset  
327 [21, 22], which is currently the largest publicly available  
328 RRG dataset, containing 377,110 chest X-ray images and  
329 227,835 radiology reports from 63,478 patients. MIMIC-  
330 CXR provides 14 uncertainty labels extracted by CheXpert

Table 1. The performance improvement of SOTA RRG methods by integrating RRGCL on the MIMIC-CXR dataset.

Method	NLG Metrics					CE Metrics			
	BL-1 $\uparrow$	BL-4 $\uparrow$	MTR $\uparrow$	RG-L $\uparrow$	Avg. $\Delta$	P $\uparrow$	R $\uparrow$	F1 $\uparrow$	Avg. $\Delta$
R <sup>2</sup> GEN[7] (EMNLP’2020)	0.353	0.103	0.142	0.277	-	0.333	0.273	0.276	-
R <sup>2</sup> GEN + RRGCL	0.361	0.108	0.152	0.284	<b>+3.429%</b>	0.341	0.280	0.282	<b>+2.380%</b>
KiUT[16] (CVPR’2023)	0.393	0.113	0.160	0.285	-	0.371	0.318	0.321	-
KiUT + RRGCL	0.402	0.117	0.169	0.293	<b>+3.260%</b>	0.382	0.327	0.329	<b>+2.772%</b>
RGRG[47] (CVPR’2023)	0.373	0.126	0.168	0.264	-	0.495	0.475	0.447	-
RGRG + RRGCL	0.388	0.129	0.177	0.278	<b>+4.404%</b>	0.512	0.489	0.459	<b>+4.035%</b>
COMG[12] (WACV’2024)	0.363	0.124	0.128	0.290	-	0.424	0.291	0.345	-
COMG + RRGCL	0.373	0.126	0.141	0.298	<b>+3.646%</b>	0.437	0.313	0.356	<b>+4.340%</b>
EKAGen[5] (CVPR’2024)	0.419	0.119	0.157	0.264	-	0.517	0.483	0.499	-
EKAGen + RRGCL	0.427	0.122	0.163	0.274	<b>+2.815%</b>	0.522	0.498	0.513	<b>+2.268%</b>
DACG[24] (MedIA’2025)	0.398	0.117	0.162	0.290	-	0.422	0.405	0.389	-
DACG + RRGCL	0.409	0.120	0.166	0.296	<b>+2.482%</b>	0.436	0.418	0.403	<b>+3.372%</b>

Table 2. Ablation studies on the learning difficulty measurer and sample scheduler conducted on the MIMIC-CXR dataset. The best results are bolded.

Method	Difficulty Measurer		Scheduler	NLG Metrics				CE Metrics		
	VNS-Gap	ITAS		BL-1	BL-4 $\uparrow$	MTR $\uparrow$	RG-L $\uparrow$	P $\uparrow$	R $\uparrow$	F1 $\uparrow$
Baseline (RGRG [47])	-	-	-	0.373	0.126	0.168	0.264	0.495	0.475	0.447
(a)	✓	-	Ours	0.381	0.127	0.173	0.272	0.508	0.482	0.454
(b)	-	✓	Ours	0.378	0.127	0.171	0.269	0.504	0.480	0.451
(c)	✓	✓	Baby Step	0.375	0.126	0.170	0.267	0.499	0.478	0.450
(d)	✓	✓	Linear	0.376	0.126	0.172	0.270	0.501	0.479	0.452
(e)	✓	✓	One Pass	0.375	0.126	0.169	0.267	0.500	0.478	0.449
(f)	✓	✓	Root	0.382	0.128	0.174	0.274	0.508	0.485	0.455
(g)	✓	✓	Geom	0.379	0.127	0.171	0.272	0.504	0.481	0.452
RGRG + RRGCL	✓	✓	Ours	<b>0.388</b>	<b>0.129</b>	<b>0.177</b>	<b>0.278</b>	<b>0.512</b>	<b>0.489</b>	<b>0.459</b>

[17]. We adopt 13 keywords excluding "Support Devices" and treat the keywords labeled as "1" (positive) as the labels for the corresponding samples. The construction of RRGCL and the training of the RRG models are independent. During the training phase of the learning difficulty measurer, all images are resized to  $256 \times 256$ , and the "Findings" and "Impressions" sections in the reports are stitched together. During the training phase of the RRG model, all data pre-processing follows the combined RRG methods.

## 4.2. Evaluation Metrics

Natural language generation (NLG) metrics are commonly used standards for evaluating the quality of generated reports, including BLEU-n (BL-1, BL-4) [39], METEOR (MTR) [2] and ROUGE-L (RG-L) [30]. Some studies [33, 60] point out that NLG metrics may not be well-suited for evaluating radiology report generation model. Therefore, to more accurately evaluate the proposed method, we used several clinical efficacy (CE) metrics, including preci-

sion (P), recall (R) and F1-score (F1).

## 4.3. Implementation Details

For the VNS-Gap measuring network, we use a ResNet101 [14] fine-tuned on CheXpert [17] to extract visual features. For the ITAS computation network, following Liao et al. [29], we employ a pre-trained 5-layer ResNet [14] to extract visual features. All textual features are extracted using the pre-trained Clinical BERT [1]. The hyperparameters  $\alpha$  and  $\beta$  in Eq.10 are set to 0.6 and 0.4, respectively. During the training of the RRG models, the configuration of training parameters follows the original settings of the combined RRG methods. All training processes are performed on an NVIDIA GeForce RTX A6000 GPU.

## 4.4. Improvements to SOTA RRG Models

To evaluate how RRGCL helps improve the performance of RRG models, we embed RRGCL into several SOTA RRG methods including R<sup>2</sup>GEN[7], KiUT [16], RGRG

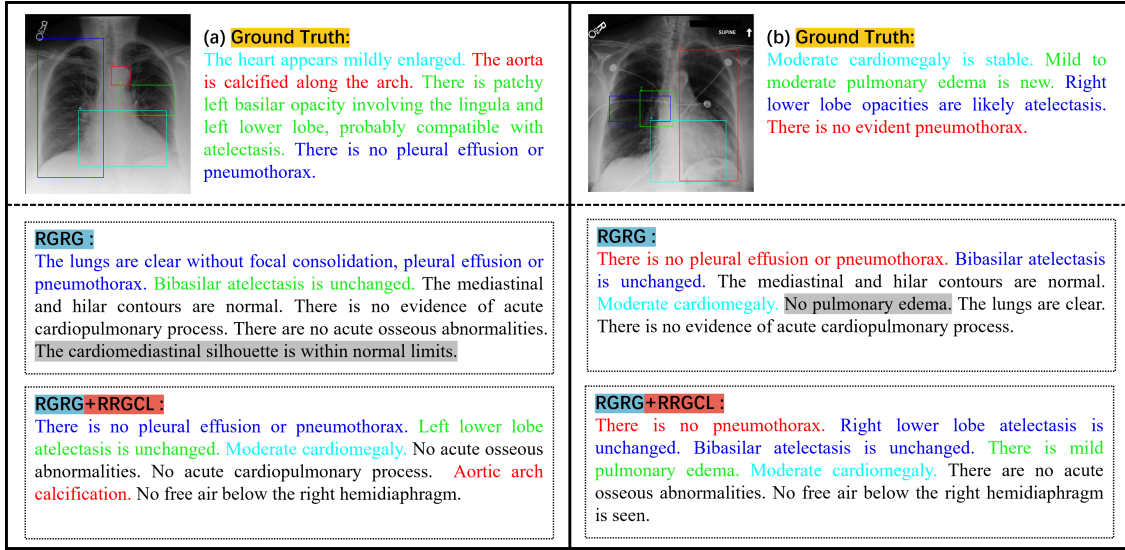


Figure 4. Two examples of reports generated by RGRG [47] and RGRG + RRGCL. The bounding boxes are provided by the ImaGenome dataset [54], and they correspond to the descriptions in the reports based on color. Sentences with gray shading represent incorrect descriptions.

Table 3. Hyperparameter sensitivity experiments on  $\alpha$  and  $\beta$  based on RGRG[47]. The best results are bolded.

$\alpha$	$\beta$	BL-1 $\uparrow$	BL-4 $\uparrow$	F1 $\uparrow$
0.2	0.8	0.382	0.128	0.451
0.4	0.6	0.385	0.128	0.454
0.6	0.4	<b>0.388</b>	<b>0.129</b>	<b>0.459</b>
0.8	0.2	0.386	0.129	0.457

[47], COMG [12], EKAGen [5] and DACG [24], and observe the quality of the generated reports. The experimental results are shown in Table 1. The results indicate that when combined with RRGCL, the performance of RRG models demonstrates significant improvements in both NLG and CE metrics. This proves that our proposed method can effectively enhance existing works and has the potential to be applied to more RRG models.

#### 4.5. Ablation Study

To validate the effectiveness of the proposed learning difficulty measurer and sample scheduler, we conduct ablation experiments. We select RGRG [47], which has the most significant performance improvement in Table 1 as the baseline, and the results are shown in Table 2. Among them, Baby Step [3], Linear [3], One Pass [3], Root [41] and Geometric Progression [40] are scheduling strategies proposed in previous works. The experimental results of methods (a) and (b) demonstrate that our proposed VNS-Gap and ITAS can effectively characterize the learning difficulty of RRG

data, with VNS-Gap being relatively more efficient. Compared to the baseline, methods (c)-(g) all achieve some improvements, while the use of our proposed scheduler yield the best performance. It is worth noting that our proposed scheduling function further allocates class-level scheduling functions based on the Root function. Compared to method (f), our method achieve improved performance because the proposed scheduling strategy not only ensures a smooth increase in learning difficulty but also maintains a balanced distribution of various class in each training epoch.

#### 4.6. Hyperparameter Sensitivity Experiment

To determine the optimal settings for the hyperparameters  $\alpha$  and  $\beta$  in Eq.10, we conduct a hyperparameter sensitivity experiment based on RGRG [47]. The results are shown in Table 3. It can be observed that the model achieves the best performance when  $\alpha$  and  $\beta$  are set to 0.6 and 0.4, respectively.

#### 4.7. Qualitative Analysis

To qualitatively analyze the effectiveness of RRGCL, we randomly selected two test samples and compared the reports generated by RGRG [47] and RGRG + RRGCL, the results are shown in Figure 4. In example (a), RGRG fails to detect ‘‘aortic calcification’’ and incorrectly considers the heart size to be normal. However, after introducing RRGCL, the generated report describes both. In the image of example (b), there is a region of ‘‘pulmonary edema’’, and the report generated by RGRG + RRGCL describes it, while RGRG provides an incorrect description. The qualitative experimental results visually demonstrate the effect of

414 RRGCL in enhancing the performance of the RRG model.

## 415 5. Conclusion

416 In this work, we focus on a previously overlooked issue  
417 in radiology report generation (RRG): the impact of data  
418 imbalance in RRG training on model performance. We  
419 discover that traditional class imbalance based on sample  
420 quantity fails to adequately characterize the imbalance in  
421 RRG data. Therefore, from a curriculum learning perspec-  
422 tive, we redefine the imbalance in RRG data by learning  
423 difficulty. Our proposed measurer calculates the sample  
424 learning difficulty from the perspectives of visual neighbor-  
425 hood semantic gap and image-text alignment scores. To ad-  
426 dress the imbalance we proposed, we further introduce a  
427 class-level dynamic sample scheduler. Experimental results  
428 demonstrate that our proposed method can effectively en-  
429 hance the performance of existing RRG models, offering a  
430 new direction for RRG studies. In the future, with the re-  
431 lease of labeled RRG datasets, we will validate our method  
432 on more datasets.

## 433 References

- 434 [1] Emily Alsentzer, John Murphy, William Boag, Wei-Hung  
435 Weng, Di Jindi, Tristan Naumann, and Matthew McDermott.  
436 Publicly available clinical bert embeddings. In *Proceedings*  
437 *of the 2nd clinical natural language processing workshop*,  
438 pages 72–78, 2019. 6
- 439 [2] Satanjeev Banerjee and Alon Lavie. Meteor: An automatic  
440 metric for mt evaluation with improved correlation with hu-  
441 man judgments. In *Proceedings of the acl workshop on in-*  
442 *trinsic and extrinsic evaluation measures for machine trans-*  
443 *lation and/or summarization*, pages 65–72, 2005. 6
- 444 [3] Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Ja-  
445 son Weston. Curriculum learning. In *Proceedings of the 26th*  
446 *annual international conference on machine learning*, pages  
447 41–48, 2009. 2, 3, 7
- 448 [4] RJM Bruls and RM Kwee. Workload for radiologists during  
449 on-call hours: dramatic increase in the past 15 years. *Insights*  
450 *into imaging*, 11:1–7, 2020. 1
- 451 [5] Shenshen Bu, Taiji Li, Yuedong Yang, and Zhiming Dai.  
452 Instance-level expert knowledge and aggregate discrimina-  
453 tive attention for radiology report generation. In *Proceed-*  
454 *ings of the IEEE/CVF Conference on Computer Vision and*  
455 *Pattern Recognition*, pages 14194–14204, 2024. 2, 6, 7
- 456 [6] Rohan Chaudhury, Maria Teleki, Xiangjue Dong, and James  
457 Caverlee. Dacl: Disfluency augmented curriculum learn-  
458 ing for fluent text generation. In *Proceedings of the 2024*  
459 *Joint International Conference on Computational Linguistics,*  
460 *Language Resources and Evaluation (LREC-COLING*  
461 *2024)*, pages 4311–4321, 2024. 3
- 462 [7] Zhihong Chen, Yan Song, Tsung-Hui Chang, and Xiang  
463 Wan. Generating radiology reports via memory-driven trans-  
464 former. In *Proceedings of the 2020 Conference on Empirical*  
465 *Methods in Natural Language Processing (EMNLP)*, pages  
466 1439–1449, 2020. 2, 6
- [8] Francesco Dalla Serra, Chaoyang Wang, Fani Deligianni,  
Jeff Dalton, and Alison Q O’Neil. Controllable chest x-  
ray report generation from longitudinal representations. In  
*The 2023 Conference on Empirical Methods in Natural Lan-*  
*guage Processing*, 2023. 2
- [9] Alexey Dosovitskiy. An image is worth 16x16 words:  
Transformers for image recognition at scale. *arXiv preprint*  
*arXiv:2010.11929*, 2020. 2
- [10] Mark Endo, Rayan Krishnan, Viswesh Krishna, Andrew Y  
Ng, and Pranav Rajpurkar. Retrieval-based chest x-ray report  
generation using a pre-trained contrastive language-image  
model. In *Machine Learning for Health*, pages 209–219.  
PMLR, 2021. 2
- [11] Xiang Geng, Yu Zhang, Jiahuan Li, Shujian Huang, Hao  
Yang, Shimin Tao, Yimeng Chen, Ning Xie, and Jiajun Chen.  
Denosing pre-training for machine translation quality esti-  
mation with curriculum learning. In *Proceedings of the AAAI*  
*Conference on Artificial Intelligence*, pages 12827–12835,  
2023. 3
- [12] Tiancheng Gu, Dongnan Liu, Zhiyuan Li, and Weidong Cai.  
Complex organ mask guided radiology report generation. In  
*Proceedings of the IEEE/CVF winter conference on appli-*  
*cations of computer vision*, pages 7995–8004, 2024. 2, 6,  
7
- [13] Tiancheng Gu, Kaicheng Yang, Xiang An, Ziyong Feng,  
Dongnan Lin, and Weidong Cai. Orid: Organ-regional infor-  
mation driven framework for radiology report generation. In  
*2025 IEEE/CVF Winter Conference on Applications of Com-*  
*puter Vision (WACV)*, pages 378–387. IEEE, 2025. 2
- [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun.  
Deep residual learning for image recognition. In *Proceed-*  
*ings of the IEEE conference on computer vision and pattern*  
*recognition*, pages 770–778, 2016. 6
- [15] Yuge Huang, Yuhan Wang, Ying Tai, Xiaoming Liu,  
Pengcheng Shen, Shaoxin Li, Jilin Li, and Feiyue Huang.  
Curricularface: adaptive curriculum learning loss for deep  
face recognition. In *proceedings of the IEEE/CVF con-*  
*ference on computer vision and pattern recognition*, pages  
5901–5910, 2020. 3
- [16] Zhongzhen Huang, Xiaofan Zhang, and Shaoting Zhang.  
Kiut: Knowledge-injected u-transformer for radiology report  
generation. In *Proceedings of the IEEE/CVF Conference*  
*on Computer Vision and Pattern Recognition*, pages 19809–  
19818, 2023. 2, 6
- [17] Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, Sil-  
viana Ciurea-Ilcus, Chris Chute, Henrik Marklund, Behzad  
Haghighi, Robyn Ball, Katie Shpanskaya, et al. Chexpert:  
A large chest radiograph dataset with uncertainty labels and  
expert comparison. In *Proceedings of the AAAI conference*  
*on artificial intelligence*, pages 590–597, 2019. 6
- [18] Mohamed Ishmael Belghazi, Aristide Baratin, Sai Rajeswar,  
Sherjil Ozair, Yoshua Bengio, Aaron Courville, and R De-  
von Hjelm. Mine: mutual information neural estimation.  
*arXiv e-prints*, pages arXiv–1801, 2018. 4
- [19] Jaehwan Jeong, Katherine Tian, Andrew Li, Sina Hartung,  
Subathra Adithan, Fardad Behzadi, Juan Calle, David Os-  
ayande, Michael Pohlen, and Pranav Rajpurkar. Multimodal

- image-text matching improves retrieval-based chest x-ray report generation. In *Medical Imaging with Deep Learning*, pages 978–990. PMLR, 2024. 2
- [20] Amelia Jiménez-Sánchez, Diana Mateus, Sonja Kirchoff, Chlodwig Kirchoff, Peter Biberthaler, Nassir Navab, Miguel A González Ballester, and Gemma Piella. Medical-based deep curriculum learning for improved fracture classification. In *Medical Image Computing and Computer Assisted Intervention–MICCAI 2019: 22nd International Conference, Shenzhen, China, October 13–17, 2019, Proceedings, Part VI 22*, pages 694–702. Springer, 2019. 3
- [21] Alistair EW Johnson, Tom J Pollard, Seth J Berkowitz, Nathaniel R Greenbaum, Matthew P Lungren, Chih-ying Deng, Roger G Mark, and Steven Horng. Mimic-cxr, a de-identified publicly available database of chest radiographs with free-text reports. *Scientific data*, 6(1):317, 2019. 2, 5
- [22] Alistair EW Johnson, Tom J Pollard, Nathaniel R Greenbaum, Matthew P Lungren, Chih-ying Deng, Yifan Peng, Zhiyong Lu, Roger G Mark, Seth J Berkowitz, and Steven Horng. Mimic-cxr-jpg, a large publicly available database of labeled chest radiographs. *arXiv preprint arXiv:1901.07042*, 2019. 2, 5
- [23] Ömer Kasalak, Haider Alnahwi, Romy Toxopeus, Jan P Pennings, Derya Yakar, and Thomas C Kwee. Work overload and diagnostic errors in radiology. *European Journal of Radiology*, 167:111032, 2023. 1
- [24] Wangyu Lang, Zhi Liu, and Yijia Zhang. Dacg: Dual attention and context guidance model for radiology report generation. *Medical Image Analysis*, 99:103377, 2025. 2, 6, 7
- [25] Nana Lange, Flavius Frasinca, and Maria Mihaela Truşcă. Curriculum learning for a hybrid approach for aspect-based sentiment analysis. *Expert Systems with Applications*, page 129669, 2025. 3
- [26] Chaoyi Li, Meng Li, Can Peng, and Brian C Lovell. Dynamic curriculum learning via in-domain uncertainty for medical image classification. In *International conference on medical image computing and computer-assisted intervention*, pages 747–757. Springer, 2023. 3
- [27] Mingjie Li, Haokun Lin, Liang Qiu, Xiaodan Liang, Ling Chen, Abdulmotaleb Elsaddik, and Xiaojun Chang. Contrastive learning with counterfactual explanations for radiology report generation. In *European Conference on Computer Vision*, pages 162–180. Springer, 2024. 2
- [28] Shan hao Li, Bang Yang, and Yuexian Zou. Adaptive curriculum learning for video captioning. *IEEE Access*, 10:31751–31759, 2022. 3
- [29] Ruizhi Liao, Daniel Moyer, Miriam Cha, Keegan Quigley, Seth Berkowitz, Steven Horng, Polina Golland, and William M Wells. Multimodal representation learning via maximization of local mutual information. In *Medical Image Computing and Computer Assisted Intervention–MICCAI 2021: 24th International Conference, Strasbourg, France, September 27–October 1, 2021, Proceedings, Part II 24*, pages 273–283. Springer, 2021. 4, 6
- [30] Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81, 2004. 6
- [31] Aohan Liu, Yuchen Guo, Jun-hai Yong, and Feng Xu. Multi-grained radiology report generation with sentence-level image-language contrastive learning. *IEEE Transactions on Medical Imaging*, 43(7):2657–2669, 2024. 2
- [32] Fenglin Liu, Shen Ge, Yuexian Zou, and Xian Wu. Competence-based multimodal curriculum learning for medical report generation. *arXiv preprint arXiv:2206.14579*, 2022. 3
- [33] Guanxiong Liu, Tzu-Ming Harry Hsu, Matthew McDermott, Willie Boag, Wei-Hung Weng, Peter Szolovits, and Marzyeh Ghassemi. Clinically accurate chest x-ray report generation. In *Machine Learning for Healthcare Conference*, pages 249–269. PMLR, 2019. 6
- [34] Kang Liu, Zhuoqi Ma, Xiaolu Kang, Yunan Li, Kun Xie, Zhicheng Jiao, and Qiguang Miao. Enhanced contrastive learning with multi-view longitudinal data for chest x-ray report generation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 10348–10359, 2025. 2
- [35] Chongwen Lyu, Chengjian Qiu, Kai Han, Saisai Li, Victor S Sheng, Huan Rong, Yuqing Song, Yi Liu, and Zhe Liu. Automatic medical report generation combining contrastive learning and feature difference. *Knowledge-Based Systems*, 305:112630, 2024. 2
- [36] Pablo Messina, Pablo Pino, Denis Parra, Alvaro Soto, Cecilia Besa, Sergio Uribe, Marcelo Andía, Cristian Tejos, Claudia Prieto, and Daniel Capurro. A survey on deep learning and explainability for automatic report generation from medical images. *ACM Computing Surveys (CSUR)*, 54(10s):1–40, 2022. 2
- [37] Hoang Nguyen, Dong Nie, Taivanbat Badamdorj, Yujie Liu, Yingying Zhu, Jason Truong, and Li Cheng. Automated generation of accurate & fluent medical x-ray reports. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3552–3569, 2021. 2
- [38] Toru Nishino, Yasuhide Miura, Tomoki Taniguchi, Tomoko Ohkuma, Yuki Suzuki, Shoji Kido, and Noriyuki Tomiyama. Factual accuracy is not enough: Planning consistent description order for radiology report generation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 7123–7138, 2022. 2
- [39] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318, 2002. 6
- [40] Gustavo Penha and Claudia Hauff. Curriculum learning strategies for ir: An empirical study on conversation response ranking. In *European conference on information retrieval*, pages 699–713. Springer, 2020. 7
- [41] Emmanouil Antonios Platanios, Otilia Stretcu, Graham Neubig, Barnabas Poczos, and Tom Mitchell. Competence-based curriculum learning for neural machine translation. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, pages 1162–1172, 2019. 7

- 639 [42] Emmanouil Antonios Platanios, Otilia Stretcu, Graham Neu- 697  
640 big, Barnabas Poczos, and Tom M Mitchell. Competence- 698  
641 based curriculum learning for neural machine translation. 699  
642 *arXiv preprint arXiv:1903.09848*, 2019. 5 700
- 643 [43] Gabriel Reale-Nosei, Elvira Amador-Domínguez, and 701  
644 Emilio Serrano. From vision to text: A comprehensive re- 702  
645 view of natural image captioning in medical diagnosis and 703  
646 radiology report generation. *Medical Image Analysis*, page 704  
647 103264, 2024. 2 705
- 648 [44] Phillip Sloan, Philip Clatworthy, Edwin Simpson, and Majid 706  
649 Mirmehdi. Automated radiology report generation: A review 707  
650 of recent advances. *IEEE Reviews in Biomedical Engineer- 708*  
651 *ing*, 2024. 2 709
- 652 [45] Petru Soviany, Radu Tudor Ionescu, Paolo Rota, and Nicu 710  
653 Sebe. Curriculum learning: A survey. *International Journal 711*  
654 *of Computer Vision*, 130(6):1526–1565, 2022. 3 712
- 655 [46] Liwen Sun, James Jialun Zhao, Wenjing Han, and Chenyan 713  
656 Xiong. Fact-aware multimodal retrieval augmentation for 714  
657 accurate medical radiology report generation. In *Proceed- 715*  
658 *ings of the 2025 Conference of the Nations of the Americas 716*  
659 *Chapter of the Association for Computational Linguistics: 717*  
660 *Human Language Technologies (Volume 1: Long Papers)*, 718  
661 pages 643–655, 2025. 2 719
- 662 [47] Tim Tanida, Philip Müller, Georgios Kaissis, and Daniel 720  
663 Rueckert. Interactive and explainable region-guided radiol- 721  
664 ogy report generation. In *Proceedings of the IEEE/CVF Con- 722*  
665 *ference on Computer Vision and Pattern Recognition*, pages 723  
666 7433–7442, 2023. 2, 6, 7 724
- 667 [48] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszko- 725  
668 reit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia 726  
669 Polosukhin. Attention is all you need. *Advances in neural 727*  
670 *information processing systems*, 30, 2017. 2 728
- 671 [49] Jun Wang, Abhir Bhalerao, and Yulan He. Cross-modal pro- 729  
672 totype driven network for radiology report generation. In 730  
673 *European Conference on Computer Vision*, pages 563–579. 731  
674 Springer, 2022. 2 732
- 675 [50] Xin Wang, Yudong Chen, and Wenwu Zhu. A survey on 733  
676 curriculum learning. *IEEE transactions on pattern analysis 734*  
677 *and machine intelligence*, 44(9):4555–4576, 2021. 3 735
- 678 [51] Xinyi Wang, Graziela Figueredo, Ruizhe Li, Wei Emma 736  
679 Zhang, Weitong Chen, and Xin Chen. A survey of deep- 737  
680 learning-based radiology report generation using multimodal 738  
681 inputs. *Medical Image Analysis*, page 103627, 2025. 2 739
- 682 [52] Zhanyu Wang, Luping Zhou, Lei Wang, and Xiu Li. A 740  
683 self-boosting framework for automated radiographic report 741  
684 generation. In *Proceedings of the IEEE/CVF Conference 742*  
685 *on Computer Vision and Pattern Recognition*, pages 2433– 743  
686 2442, 2021. 2 734
- 687 [53] Zhanyu Wang, Lingqiao Liu, Lei Wang, and Luping Zhou. 744  
688 Metransformer: Radiology report generation by transformer 745  
689 with multiple learnable expert tokens. In *Proceedings of 746*  
690 *the IEEE/CVF Conference on Computer Vision and Pattern 747*  
691 *Recognition*, pages 11558–11567, 2023. 2 748
- 692 [54] Joy T Wu, Nkechinyere N Agu, Ismini Lourentzou, Arjun 749  
693 Sharma, Joseph A Paguio, Jasper S Yao, Edward C Dee, 750  
694 William Mitchell, Satyananda Kashyap, Andrea Giovannini, 751  
695 et al. Chest imagenome dataset for clinical reasoning. *arXiv 752*  
696 *preprint arXiv:2108.00316*, 2021. 7 753
- [55] Linhui Xiao, Xiaoshan Yang, Fang Peng, Ming Yan, Yaowei 697  
Wang, and Changsheng Xu. Clip-vg: Self-paced curriculum 698  
adapting of clip for visual grounding. *IEEE Transactions on 699*  
*Multimedia*, 2023. 3 700
- [56] Ting Xiao, Lei Shi, Peng Liu, Zhe Wang, and Chenjia Bai. 701  
Radiology report generation via multi-objective preference 702  
optimization. In *Proceedings of the AAAI Conference on Ar- 703*  
*tificial Intelligence*, pages 8664–8672, 2025. 2 704
- [57] Lin Yang, Yi Shen, Yue Mao, and Longjun Cai. Hybrid 705  
curriculum learning for emotion recognition in conversation. 706  
In *Proceedings of the AAAI Conference on Artificial Intelli- 707*  
*gence*, pages 11595–11603, 2022. 3 708
- [58] Xingyi Yang, Muchao Ye, Quanzeng You, and Fenglong Ma. 709  
Writing by memorizing: Hierarchical retrieval-based med- 710  
ical report generation. In *Proceedings of the 59th Annual 711*  
*Meeting of the Association for Computational Linguistics 712*  
*and the 11th International Joint Conference on Natural Lan- 713*  
*guage Processing (Volume 1: Long Papers)*, pages 5000– 714  
5009, 2021. 2 715
- [59] Xinyu Yang, Tilo Burghardt, and Majid Mirmehdi. Dynamic 716  
curriculum learning for great ape detection in the wild. *Inter- 717*  
*national Journal of Computer Vision*, 131(5):1163–1181, 718  
2023. 3 719
- [60] Feiyang Yu, Mark Endo, Rayan Krishnan, Ian Pan, Andy 720  
Tsai, Eduardo Pontes Reis, Eduardo Kaiser Ururahy Nunes 721  
Fonseca, Henrique Min Ho Lee, Zahra Shakeri Hossein 722  
Abad, Andrew Y Ng, et al. Evaluating progress in automatic 723  
chest x-ray radiology report generation. *Patterns*, 4(9), 2023. 724  
6 725
- [61] Yixiao Zhang, Xiaosong Wang, Ziyue Xu, Qihang Yu, Alan 726  
Yuille, and Daguang Xu. When radiology report generation 727  
meets knowledge graph. In *Proceedings of the AAAI con- 728*  
*ference on artificial intelligence*, pages 12910–12917, 2020. 729  
2 730
- [62] Qin Zhou, Guoyan Liang, Xindi Li, Jingyuan Chen, Zhe 731  
Wang, Chang Yao, and Sai Wu. Learnable retrieval enhanced 732  
visual-text alignment and fusion for radiology report genera- 733  
tion. In *Proceedings of the IEEE/CVF International Confer- 734*  
*ence on Computer Vision*, pages 22529–22538, 2025. 2 735